

## Overview

### Visual Grounding of phrases:

Localize any textual query into a given image.



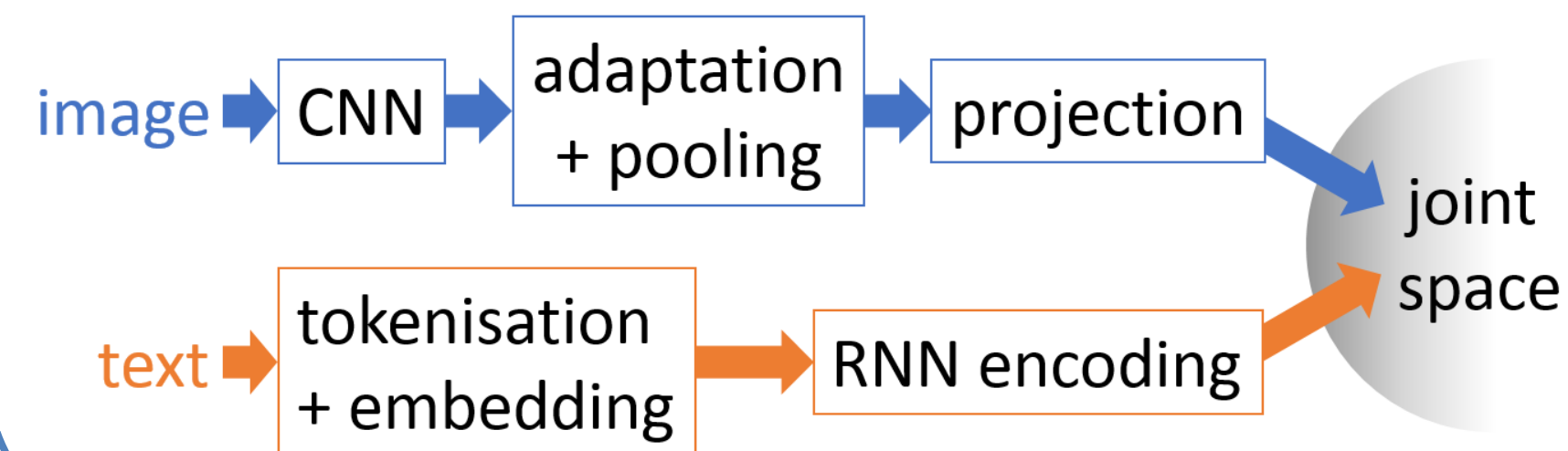
### Cross-modal retrieval:

Query: A cat on a sofa



### Approach:

- Learning image and text joint embedding space.
- Visual grounding relying on the spatial-textual information modeling.
- Cross-modal retrieval leveraging the semantic space and the visual and textual alignment.



## References

- [1] K. He et al. Deep residual learning for image recognition. CVPR, 2016.
- [2] T. Durand et al. Weldon: Weakly supervised learning of deep convolutional neural networks. CVPR, 2016.
- [3] R. Kiros et al. Skip-thought vectors. NIPS, 2015.
- [4] T. Lei et al. Training RNNs as fast as CNNs. arXiv, 2017.
- [5] A. Eisenschat et al. Linking image and text with 2-way nets. CVPR, 2017.
- [6] F. Faghri et al. VSE++: Improved visual-semantic embeddings. arXiv, 2017.
- [7] F. Xiao et al. Weakly-supervised visual grounding of phrases with linguistic structures. CVPR, 2017.

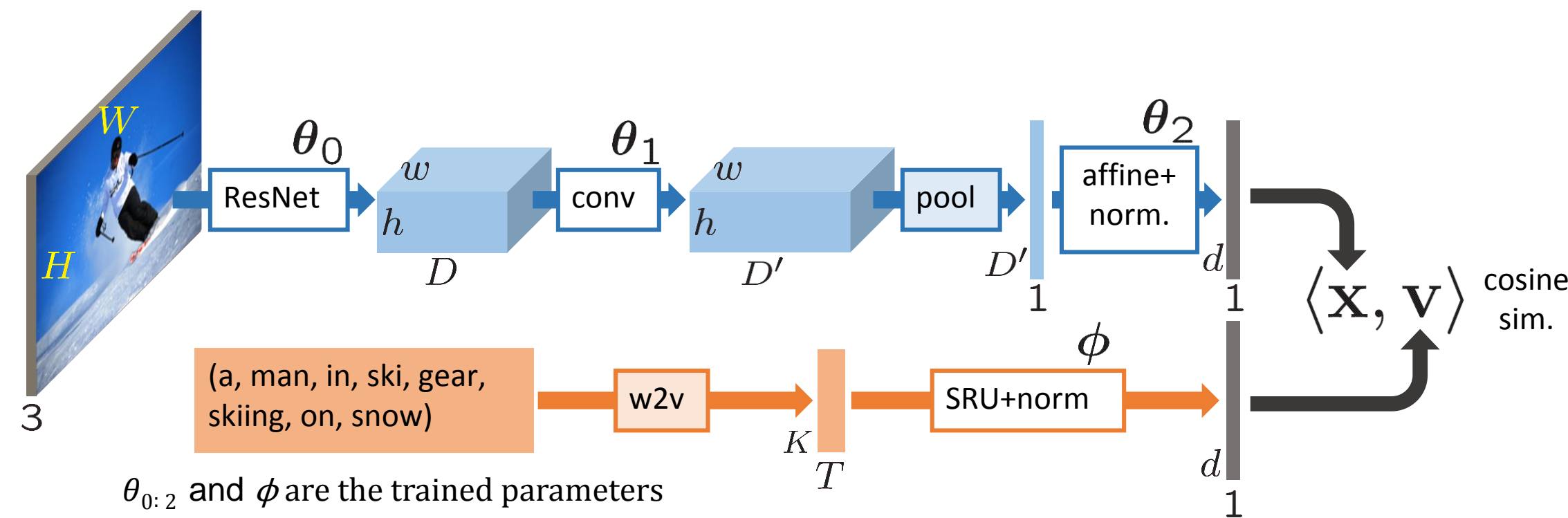
## Semantic Embedding Model

### Visual pipeline:

- ResNet-152 [1] pretrained.
- Weldon [2] spatial pooling.
- Affine projection and normalization.

### Textual pipeline:

- Pretrained word embedding [3].
- Simple Recurrent Unit (SRU) [4].
- Normalization.



Loss for a batch  $\mathcal{B} = \{(\mathbf{I}_n, \mathbf{S}_n)\}_{n \in \mathcal{B}}$  of image sentence pairs:

$$\mathcal{L}(\Theta; \mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \left( \max_{m \in C_n \cap B} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \max_{m \in D_n \cap B} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right)$$

### With:

- $\text{loss}(\mathbf{y}, \mathbf{z}, \mathbf{z}') = \max\{0, \alpha - \langle \mathbf{y}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z}' \rangle\}$
- $C_n$  (resp.  $D_n$ ) set of indices of caption (resp. image) unrelated to  $n$ -th element.

## Localization

### Visual grounding module:

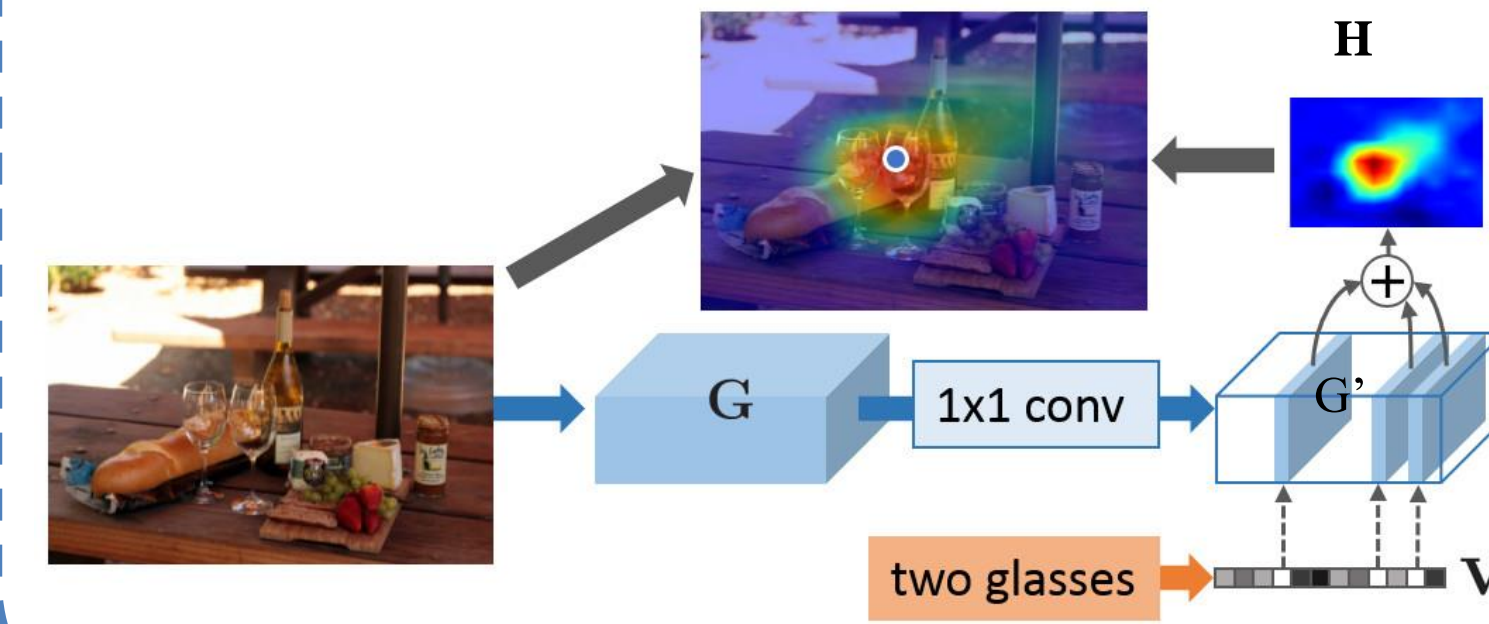
- Weakly supervised, with no additional training.
- Using the embedding space to select convolutional activation maps.

### Generation of heatmap $\mathbf{H}$ :

$$\mathbf{G}'[i, j, :] = \text{AG}[i, j, :], \forall (i, j) \in [1, w] \times [1, h]$$

$$\mathbf{H} = \sum_{u \in K(\mathbf{v})} |\mathbf{v}[u]| * \mathbf{G}'[:, :, u]$$

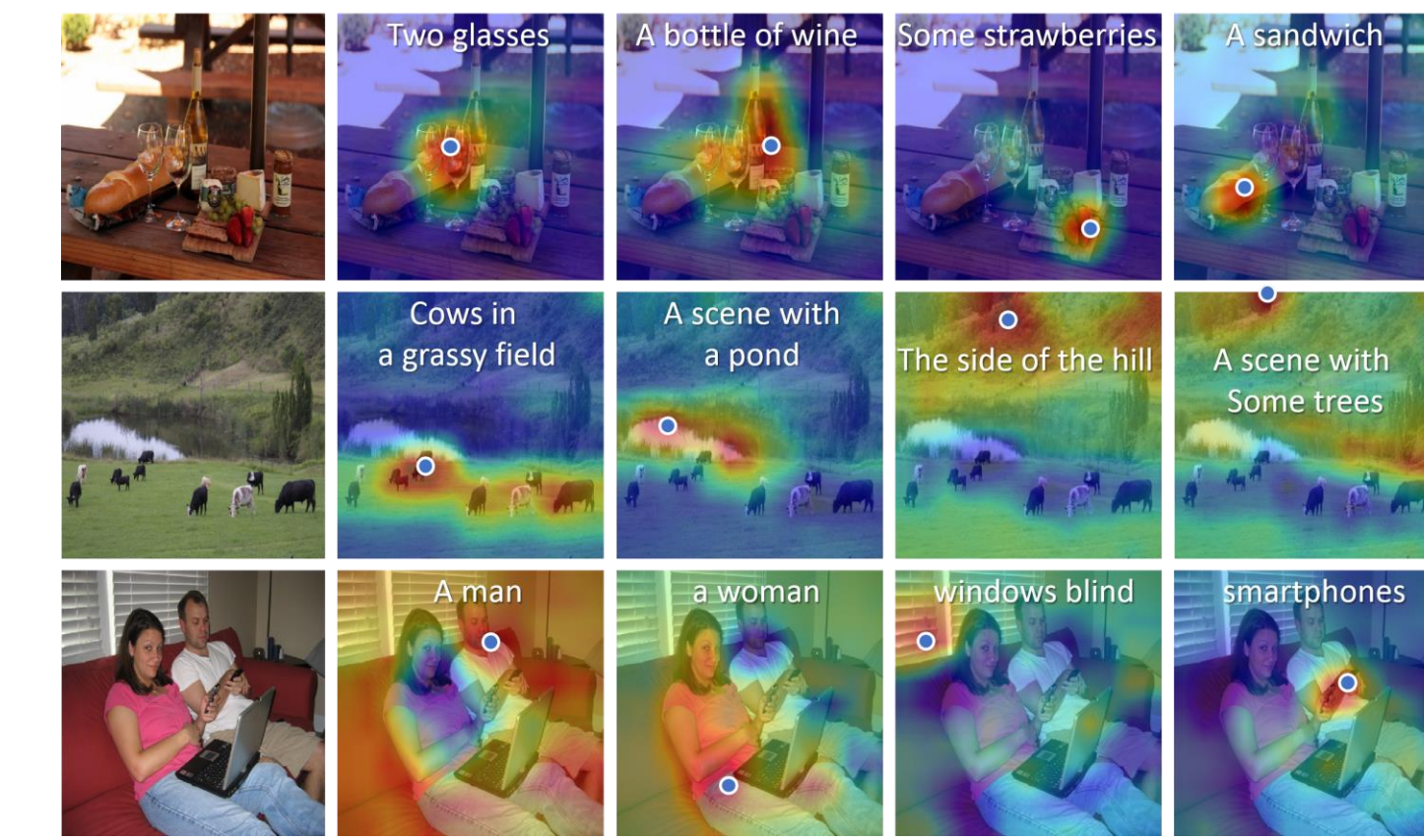
$K(\mathbf{v})$  the set of the indices of its  $k$  largest entries



## Visual results

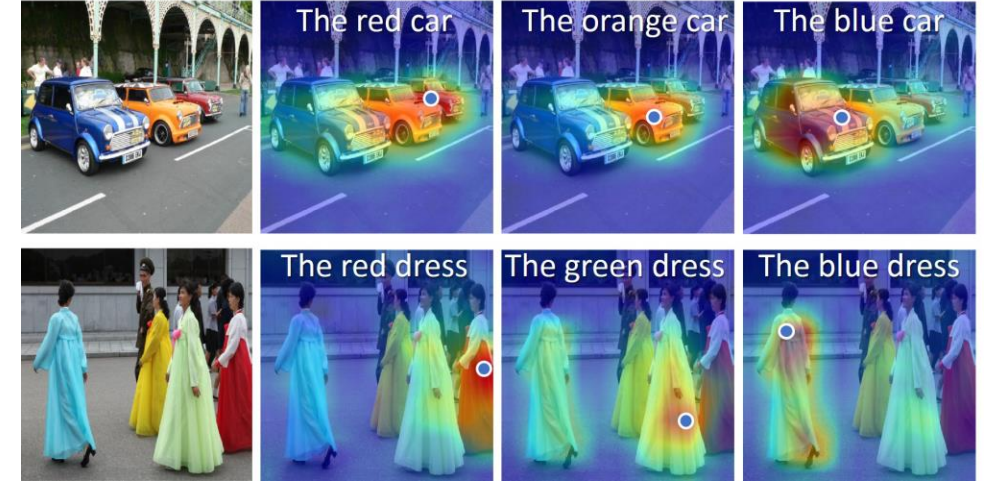
### Visual grounding examples:

- Generating multiple heat maps with different textual queries.

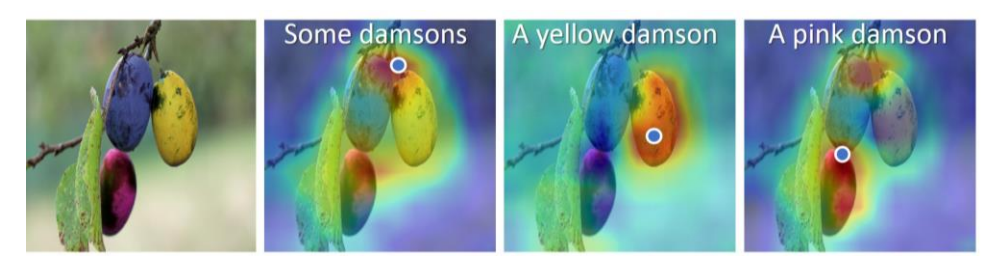


### Toward zero-shot localization:

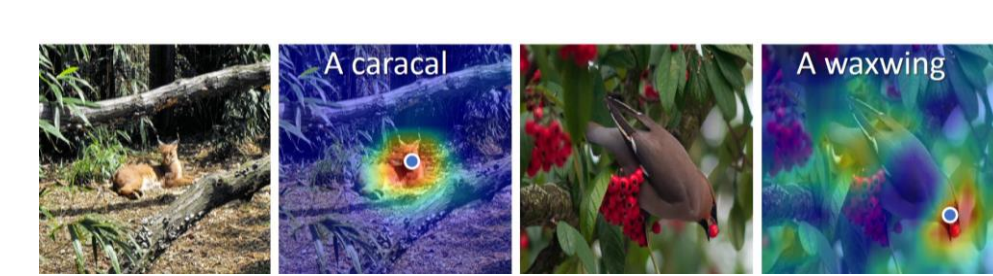
- Emergence of colors understanding:



- Even on artificial images:



- Generalization to unseen elements:



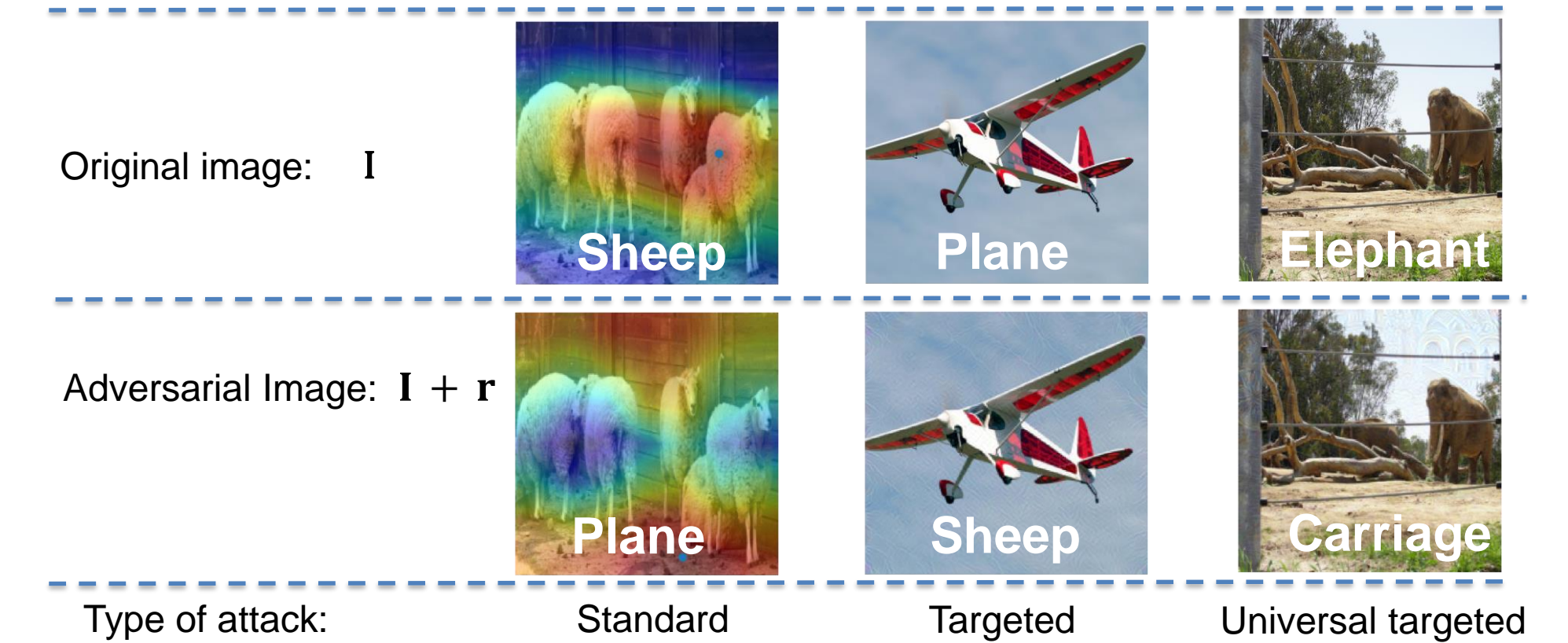
## Sensitivity to adversarial attacks

### Finding adversarial example:

Optimizing noise  $\mathbf{r}$  over input image  $\mathbf{I}$  resulting in its representation  $F(\mathbf{I} + \mathbf{r}; \theta_{0:2})$  being displaced toward  $\mathbf{y}$  in the embedding space.

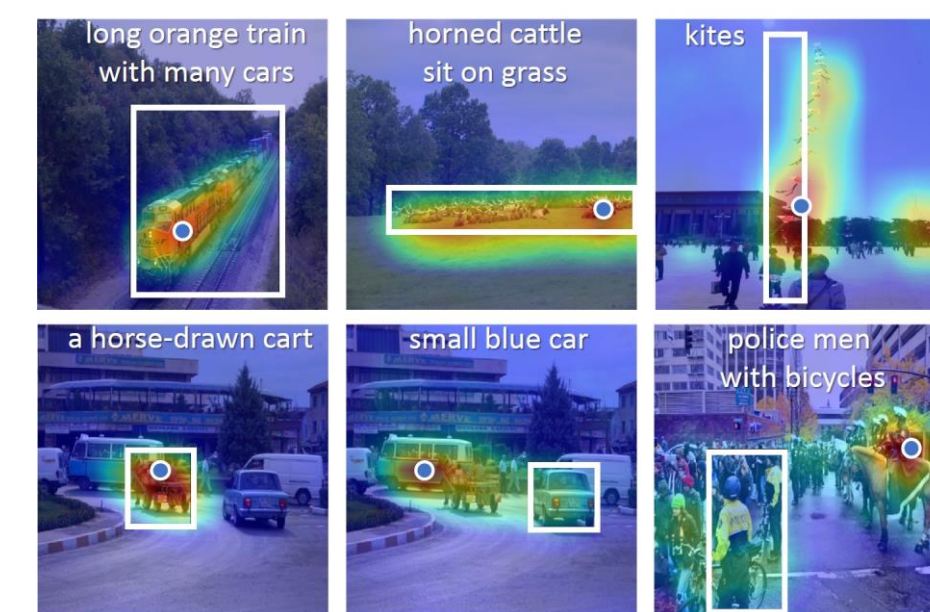
$$\max_{\mathbf{r}} \langle F(\mathbf{I} + \mathbf{r}; \theta_{0:2}), \mathbf{y} \rangle \text{ s.t. } \|\mathbf{r}\|_2 \leq \epsilon$$

### Examples of adversarial attacks

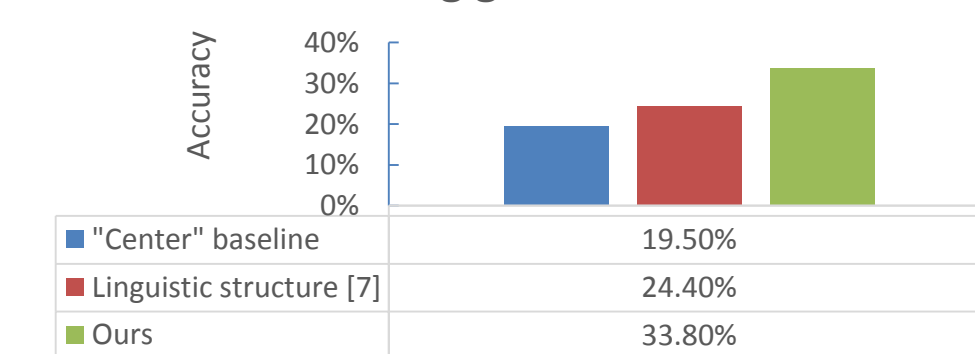


Closest word in the embedding overlaid in white.

### The pointing game: Localizing phrases corresponding to subregions of the image.

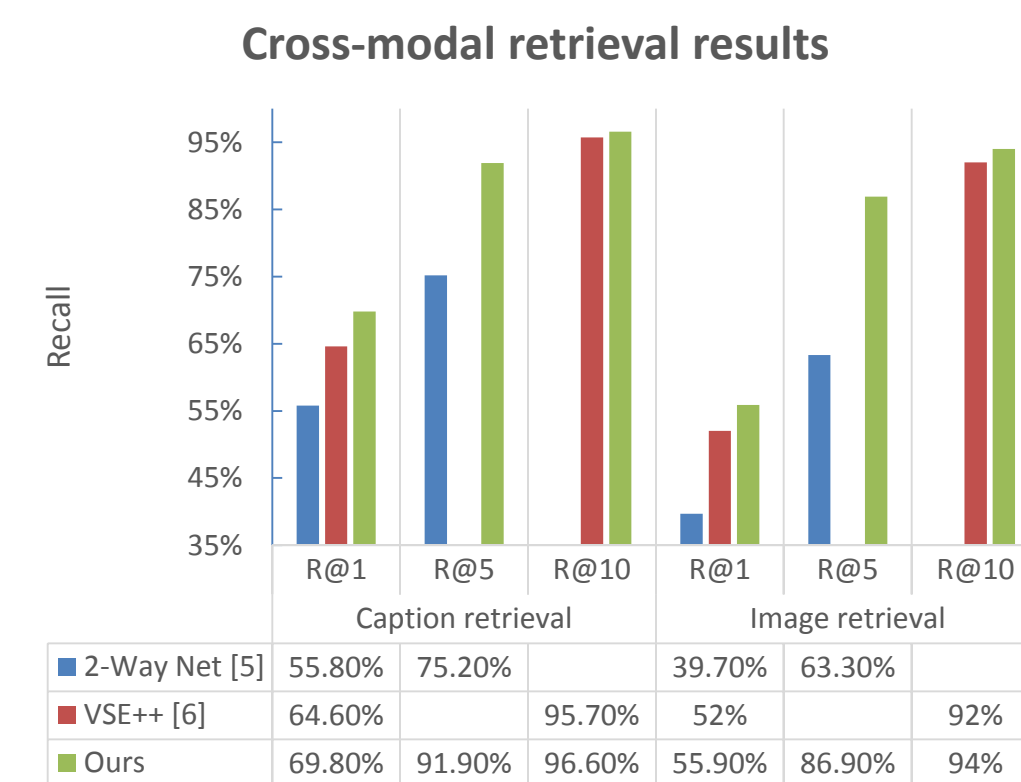


### Pointing game results



## Evaluation

### Cross-modal retrieval: Evaluated on MS-CoCo image/caption pairs.



### Performance boost coming from:

- Architecture choice: SRU and Weldon spatial pooling.
- Efficient learning strategy: hard negative loss.

### Ablation study: cross modal retrieval results

